

Exploring the repeat protein universe through computational protein design

Authors

TJ Brunette* 1,2, Fabio Parmeggiani* 1,2, Po-Ssu Huang* 1,2, Gira Bhabha 3, Damian Ekiert 4, Susan E. Tsutakawa 5, Greg L. Hura 6, John A. Tainer 5,6 and David Baker 1,2,7

*These authors contributed equally to this work

Affiliation

1 Department of Biochemistry, University of Washington, Seattle, WA 98195, USA

2 Institute for Protein Design, University of Washington, Seattle, WA 98195, USA

3 Department of Cellular and Molecular Pharmacology, UCSF, San Francisco, CA 94158, USA

4 Department of Microbiology and Immunology, UCSF, San Francisco, CA 94158, USA

5 Life Sciences Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA

6 Department of Molecular and Cellular Oncology, The University of Texas M. D. Anderson Cancer Center, 1515 Holcombe Blvd., Houston, Texas 77030, USA

7 Howard Hughes Medical Institute, University of Washington, Seattle, WA 98195, USA

Contents

Supplementary Discussion 1 | Computational protocol

Supplementary Table 1 | Detailed protocol used for each design

Supplementary Discussion 2 | Geometric parameters of Designed Helical Repeat proteins

Supplementary Table 2 | Global geometric parameters

Supplementary Table 3 | Local geometric parameters

Supplementary Discussion 3 | Structure and sequence comparison

Supplementary Table 4 | Sequence comparison

Supplementary Discussion 4 | Structure determination remarks

Supplementary Table 5 | Summary crystallization

Supplementary Tables 6-14 | Crystallographic data collection

Supplementary Discussion 5 | Small Angle X-ray Scattering analysis

Supplementary Table 15 | Summary of SAXS analysis

Supplementary Table 16 | Sequences

Supplementary References

Supplementary Discussion 1 | Computational protocol

We have developed a method for construction of Designed Helical Repeats (DHRs) depicted in Extended Data Fig. 1 and described below. We designed proteins based on repeating units formed by two helices and two loops. For all proteins this design process was completely automated and no manual refinement was involved. Using this protocol 69 proteins with diverse architectures were selected from the *in silico* candidates. For 14 models, an additional version that included disulphide bonds was selected, for a final list of 83 proteins that were experimentally tested. This design method has progressed over the duration of this research and only the final design method is described below. The status of the computational design method when each design was experimentally tested is illustrated in Supplementary Table 1.

Each of the following sections describes one step in Rosetta_examples and corresponds to the flow chart in Extended Data Fig. 1. The fraction of models passing each design stage and their distribution according to helix and loop length are shown in Extended Data Fig. 2. The Rosetta design code for each step is provided as Supplementary Material: Rosetta_examples, which requires Rosetta git version c876538.

1 Backbone Design

The backbone design stage employs a simplified side chain representation (centroid)¹. The backbone assembly procedure begins by picking fragments harvested directly from a non-redundant set of structures from PDB². The fragments contain only residues that fall into the space of phi-psi backbone angles of either helices or loops depending on the desired secondary structure. Loop fragments could be further specified to fall within desired ABEGO bins³ as described by Koga *et al.*⁴.

The fragments were assembled using a Monte-Carlo sampling procedure that was initialized with ideal-helices and extended loops. After every fragment sampling step, which was allowed only in the first repeat unit and at the junction between the first and the second units, the change was propagated to all downstream repeats and scored. The score function we used considered van der Waals interactions, packing, values of backbone dihedral angles, and radius of gyration (RG) that was applied to only the first and second repeat-unit (RG-local). The RG term promotes the formation of globular proteins so applying RG to the whole model produced only highly curved structures. The sampling procedure in the database used 1500 Monte Carlo fragment insertions and was further improved to 3200 steps ordered as following: 100 Monte Carlo moves with 9 residue fragments then 100 moves with 3 residue fragments, both allowed only in loops. The loop sampling was followed by 1500 moves with 9 residue fragments and 1500 moves with

3 residue fragments, both in helices and loops (improved sampling). The improvements resulted in a 3.3 times increase of acceptance at the centroid stage. The backbone was represented as poly-tyrosine during the centroid building, maintaining enough space within the core to accommodate both small and large side chains in the design step.

Using this procedure we designed 2.88 million backbones by making 500 structures for each of 5776 different secondary structure combination.

2 Backbone quality filter: RMSD loop threshold and motif score

Designed backbones were screened for native-like features. First, loops were checked so that there was at least one 9-residue fragment from the PDB database within 0.4Å RMSD on every position in the structure (RMSD loop threshold). To do this we used the worst9mer filter in Rosetta¹. Second, the design-ability of each residue was measured by the number of pairwise side chain interactions observed in the PDB database, considering the backbone position of the two residues involved (motif score, unpublished results). Backbones with fewer than 1.5 interactions per residue were filtered out. Of the 2.88 million initial backbones 66,776 structures passed these filters.

3 Sequence design - Fast

Starting from the filtered backbone conformations, we used one pass of Rosetta design⁵ to generate repeated sequences.

4 Packing filters – Low threshold

After completing sequence design the models were filtered out if the helices were either too far apart, creating cavities in the core (poor Rosetta holes⁶ score, > 1.75), or too close together with an alanine-rich unspecific core packing (%alanine residues > 25%). Of the 66,7776 structures that passed centroid 11,243 pass this filter.

5 Structure Profile

The structure profile biases the sequence composition towards the sequences in native proteins with similar local structure. To construct the structural profile, the sequences from the closest 100 9-residue fragments within 0.5Å RMSD to the designed structure were used. The code to construct the structural profile is included with Rosetta as generate_struct_profile.rb in tools/fragment_tools/pdb2vall. The structure profile was used in the same way as the sequence profile described by Parmeggiani *et al.*⁷

6 Sequence design – multipass

Starting from the filtered backbone conformations, we used Rosetta design to generate repeated sequences while minimizing the overall energy^{4,8}, increasing core packing as measured by Rosetta holes⁶ and improving the psipred secondary structure prediction⁹. After the first round of sequence refinement the N and C terminal repeats (capping repeats) display exposed hydrophobic residues. The sequence design procedure was-rerun for these repeats without a symmetric sequence to introduce polar amino acids.

7 Packing filters –High threshold

After completing sequence design the models were filtered out for poor packing. (holes score, < 0.5). After this stage we obtained 1980 structures.

8 Exploration of the energy landscape

The designs were validated using Rosetta *ab initio* structure prediction using Rosetta@Home^{10,11}. In Rosetta *ab initio* prediction the energy landscape is explored using independent simulations starting from an extended structure. The distribution of the simulation results is expressed in terms of energy and distance from the target fold as root mean square deviation (RMSD). A successful design produces a distribution in the shape of a funnel with the minimum corresponding to low energy and low RMSD models and no alternative minima.

For each structure, seven family members were made from the same topology, some with increased hydrogen bond potential. Proteins where multiple family members had successful simulations were selected. The member of the family with the tightest folding funnel was chosen by visual inspection and the corresponding gene was ordered for experimental testing. Extended data Fig. 3 illustrates the folding funnel and sequence diversity for one topology.

For the database we have 761 structures that have at least one family member < 3.0 RMSD from the design.

9 Add disulphides

Additional, versions with stabilizing inter-repeat disulphide bonds were also generated. Potential disulphides were scored using RosettaRemodel¹² and if the disulphide score was < 0 they were considered.

Time estimates

Backbone design: on a single core of a Xeon E5-2650 took 104.5 seconds to build a structure with a 19H-2L-20H-3L topology, the median topology in the database. With an average design time of 104.5 seconds per model, it would take 3493 compute days on a single core to generate the 2.8 million structures.

Sequence design – multipass: the multipass design of sequence and capping residues takes 2.1 hours for a model with 17 length helices and 3 length loops on a single core of a Xeon E5-2650.

Exploration of the energy landscape: on a single core of a Xeon E7-2850 @ 2.00 GHZ a model with 17 residues helices and 3 residues loops is produced in 19.7 minutes. Where the computation was run on Rosetta@Home, the average was 26.7 minutes. With 7 sequences per family and a minimum of 1000 models to suitably explore the landscape it would take 130 compute days per structure.

Supplementary Table 1 | Detailed protocol used for each design

	asymmetric topology (a)	structure profile	RMSD loop threshold	disulphide threshold < 0.0 (b)	energy landscape exploration	no csts (c)	motif score	ABEGO type (d)	improved sampling (d)
DHR 1-4,9									
DHR 5-8,10									
DHR 11-18									
DHR 19-30									
DHR 56-83									
DHR 31-55									

Dark grey cells indicate the computational method variants employed for groups of designs. The database described in section 1 of the supplementary corresponds to the technique used to make DHR56-83. (a) For DHR1-4,9,11-18 the repeat backbone at the centroid level was symmetric, with first and second helices and first and second loops having the same length and conformation. The design stage was not restricted, introducing structural and sequence variability between the two halves of the repeat. (b) A higher disulphide score threshold of 1.5 was initially used which resulted in many disulphide-containing structures being non-functional. (c) We initially used ambiguous constraints between the helices. Ambiguous constraints gave a score bonus to centroid models when a helix was within 10 Å to a helix in adjacent repeat. These constraints were found to disrupt loops and result in many structures that would not fold during simulations. (d) DHR31-55 contained a displacement between helices, which resulted in highly twisted structures. This displacement was observed when the ABEGO loop types GBB and BAB were coupled with specific helix lengths. An improved sampling strategy with increased number of Monte Carlo steps was also used in these cases. The other method variants indicated in the columns are described in supplementary section 1 (names underlined).

Supplementary Discussion 2 | Geometrical parameters of Designed Helical Repeat proteins

- 1) Global parameters
- 2) Extracting parameters from naturally occurring repeats
- 3) Local parameters

1) Global parameters

Class 3 repeat proteins, as described by Kajava A.¹³, form solenoid structures that can be described in term of global helical parameters that relate the position of one repeat to the next one: radius (r), twist or angle between adjacent repeats around the helical axis (twist, ω) and translation between adjacent repeats along the helical axis (z).

Parameters for Designed Helical Repeat proteins (DHRs) and crystal structures are indicated in Supplementary Table 2, together with the $C\alpha$ RMSD values for the complete proteins. The parameters were measured on the two central repeats using the RepeatParameter filter available in Rosetta.

Radius and twist are inversely correlated and their distribution of whole set describes a hyperbolic shape, which can be represented as two symmetric ones, when considering the handedness of the superhelix in the ω value as in Fig. 2b (+ right handed, - left handed). Handedness refers to the superhelix described by the center of mass of the repeats. z is broadly distributed, with maximum values around 16 Å (Extended Data Fig. 4).

2) Extracting parameters from naturally occurring repeats

A set of alpha-helical solenoid proteins were curated from the repeatsDB (category III.3.)¹⁴ to remove both proteins that had above 90% sequence identity^{15,16} and previously designed repeat proteins. After curation, 258 proteins remained out of 923. We then automatically extracted repeat units, which consisted of 3 subsequent repeats, that differed by less than 3 residues in length and had a high degree of structural similarity as measured by having a TM score¹⁷ of greater than 0.75. The requirement of high structural similarity cut down the number of repeat proteins to 81. Repeat units were identified by the method described by RAPHAEL¹⁸ implemented in Rosetta and improved. This method measures the distance from residues in the protein to random points placed around the protein. Equally spaced inflection points, where a residue was furthest or closest to these random points indicated the start of a repeat.

We found that inflection points occurred at random in repeat protein loops. To ensure each repeat was cut at the same location, the first residue in each repeat was chosen to be the loop-helix transition closest to the transition point. The code for this is available as

extractNativeRepeats in Rosetta after git branch c876538. After locating repeats we assigned the class name of each repeat based on the PDB assignment in the Pfam database¹⁹. The Rise/Omega/Twist parameters were calculated by superimposing the first repeat-unit onto the second using TM-align¹⁷ then calling the parameter calculators and averaging the values within the same protein. This approach does not provide an extensive coverage of all the possible curvatures for each family but an indication of the protein average values.

3) Local parameters

Local parameters describe the helix-helix interactions and, due to the repeating structures, only two interactions are needed to capture the local geometry: helix1.1-helix1.2 within a repeat and helix1.1-helix2.1 between first and second repeat. Angle between helices and distance between helix centers of mass were used as parameters, extracted with a modified version of the publicly available script <http://www.pymolwiki.org/index.php/AngleBetweenHelices>. Secondary structure definition were assigned using DSSP²⁰. Values are reported in Supplementary Table 3 for designs and crystal structures. For the two central repeats, all atoms RMSDs between crystal structures and design are reported. Repeat handedness, as defined by Kobe and Kajava²¹, indicates the rotation of the main chain going from the N- to the C-terminal around the axis connecting the repeat centers of mass.

Supplementary Discussion 3 | Structure and sequence comparison

Structural comparison of experimentally validated designs with representative repeat proteins from repeatDB¹⁴ revealed that DHRs cluster in different families than the existing repeat proteins (Extended Data Fig. 8). Additionally, designs are equally distributed between right-handed and left-handed architecture, as referred to the repeat handedness (see local parameters above), in contrast to known alpha helical repeat proteins, which are mostly right-handed. This result indicates that the handedness observed is not an intrinsic limitation of repeat proteins structures but the result of a bias during evolution.

At the sequence level, DHRs are expected to be similar to repeat proteins with high units formed by two helices, like tetratricopeptide repeats (TPR). The results of BLAST²² searches over non-redundant NR NCBI database with an e-value cutoff of 0.0001 indicate that no characterized repeat protein was identified among the top hits (Supplementary Table 4).

HHsearch²³ comparison of single repeat sequences with the whole Pfam profile database¹⁹, shows a similar trend (Supplementary Table 4), revealing that the hits are generally only very weak (Prob, E-value, P-value, score) and often rely only on a subset of position (Cols compared to repeat size). For several designs the first repeat protein hit is not among the first hits (see rank) and in some cases no similarity to repeat protein Hidden Markov Models is detected (n.d.). As reference helical proteins for this search, the ANK and TPR clans were considered, as well as spectrin and mterf families. Most helical repeat proteins (e.g. ARM, PPR, HAT, PUF) are included in the TPR clan.

Supplementary Discussion 4 | Structure determination remarks

Due to the presence of 6 cysteine residues in the native protein, the DHR5 structure was solved by sulfur single wavelength anomalous dispersion (S-SAD) using a dataset collected at 7235 eV. A search for 6 individual sulfur atoms in SHELXD gave many clear solutions that led to near complete autobuilding of a poly-alanine backbone in SHELXE, which was further elaborated using the Autobuild module of Phenix. Ultimately, the final model for DHR5 was in good agreement with the design target structure, despite our initial difficulties in phasing by molecular replacement. While the SAD data set was limited to 1.85 Å, the final model was refined against the original data set (1.25Å). Both data sets were deposited in the Protein Data Bank.

The asymmetric unit for DHR8 was found to contain 4 copies of DHR8. Although the overall structure of the 4 copies is similar, the electron density for the N-terminal helix from two of these monomers is weak, suggesting that these helices are partially disordered in the crystal. Indeed, crystal packing of these helices in the designed conformation would have led to significant steric overlap with one another. As the corresponding helices in the remaining two DHR8 monomers were well-ordered and essentially as designed, these fully ordered models were used for further analysis.

The dataset collected for DHR14 had a large non-origin Patterson peak at fractional coordinates (0.000, 0.217, 0.000), suggesting the presence of translational NCS. However, consideration of the apparent space group, unit cell parameters, and plausible solvent content strongly indicated the presence of a single copy of DHR14 in the asymmetric unit. Given the relatively low pitch of this helical design and the translational pseudosymmetry between the N- and C-terminal halves of the protein, we suspected that intramolecular pseudotranslational NCS might account for the observed Patterson peak. Ultimately, a molecular replacement solution was obtained using 4 of the 8 designed helices of DHR14, and this was sufficient to bootstrap autobuilding of the remaining backbone using SHELXE. In the final model, the helical axis of DHR14 is closely aligned with the crystallographic b axis, and pseudotranslational NCS between the N- and C-terminal repeats with a translation of ~21 Å is in good agreement with the observed fractional Patterson peak at ~0.22 along b.

Supplementary Table 5 | Summary crystallization

Protein	Concentration (mg/ml)	Cryoprotectant added	Crystallization condition	Beamline	Phasing method	Space group	resolution
DHR4	23	7.5% ethylene glycol	0.2 M Ammonium acetate, 0.1M acetate pH 4.6, 30% PEG 4000	ALS 8.3.1	Molecular replacement	P2 ₁ 2 ₁ 2 ₁	1.55
DHR5	23	17.5% ethylene glycol	1.0 M Lithium chloride, 0.1 M HEPES pH 7.0, 20% (w/v) PEG 6000	ALS 8.3.1	Molecular replacement	P2 ₁ 2 ₂ 1	1.25
DHR5	23	17.5% ethylene glycol	20% (w/v) PEG 6000, 1 M Lithium chloride, 0.1 M BisOTris pH 7.0	ALS 8.3.1	Sulfur SAD	P2 ₁ 2 ₂ 1	1.85
DHR7	19	none	0.2 M Sodium chloride, 0.1 M CHES pH 9.5, 50% (v/v) PEG 400	ALS 8.3.1	Molecular replacement	C121	2.60
DHR8	7	7.5% ethylene glycol	0.19 M Calcium chloride, 0.095 M HEPES pH 7.5, 26.6% (v/v) PEG 400, 5% (v/v) Glycerol	ALS 8.3.1	Molecular replacement	P2 ₁ 2 ₁ 2	1.80
DHR10	23	17.5% ethylene glycol	0.2 M Magnesium formate, 20% (w/v) PEG 3350	ALS 8.3.1	Molecular replacement	P2 ₁ 2 ₂ 1	1.20
DHR14	70	17.5% ethylene glycol	0.2 M Lithium nitrate, 20% (w/v) PEG 3350	ALS 8.3.1	Molecular replacement	C222 ₁	1.30
DHR18	88	none	0.1 M PhosphateOcitrate pH 4.2, 40% (v/v) PEG 300	ALS 8.3.1	Molecular replacement	P1	1.75
DHR49	31	15% ethylene glycol	0.2 M diOPotassium hydrogen phosphate, 20% (w/v) PEG 3350	ALS 8.3.1	Molecular replacement	I4 ₁	1.70
DHR53	63	12.5% ethylene glycol	26% (w/v) PEG 6000, 0.1 M Sodium citrate, pH 2.2	ALS 8.3.1	Molecular replacement	C12 ₁	1.90
DHR54	92	none	0.2 M Sodium chloride, 0.1 M Na/K phosphate pH 6.2, 50% (v/v) PEG 200	ALS 8.3.1	Molecular replacement	P2 ₁ 2 ₁ 2 ₁	1.50
DHR64	43	none	0.1 M Sodium acetate pH 4.5, 35% (v/v) MPD	ALS 8.3.1	Molecular replacement	P2	2.50
DHR64	43	25% ethylene glycol	0.1 M Sodium acetate pH 5.0, 10% (v/v) MPD	ALS 8.3.1	Molecular replacement	P6	2.90
DHR71	147	none	0.1 M CHES pH 9.5, 50% (v/v) PEG 200	ALS 8.3.1	Molecular replacement	P2 ₁ 2 ₁ 2 ₁	1.70
DHR76	66	30% ethylene glycol	0.1 M MES pH 5.0, 5% (w/v) PEG 6000	ALS 8.3.1	Molecular replacement	P2 ₁ 2 ₁ 2 ₁	3.35
DHR79	22	35% ethylene glycol	1.6 M Sodium citrate pH 6.5	ALS 8.3.1	Molecular replacement	P4 ₃ 2 ₁ 2	1.90
DHR81	39	15% ethylene glycol	0.08 M Sodium acetate pH 4.6, 20% (v/v) Glycerol	ALS 8.2.1	Molecular replacement	I23	2.05

Supplementary Table 6. Data collection and refinement statistics for DHR4 and DHR5

Data collection	DHR4	DHR5
Beamline	ALS 8.3.1	ALS 8.3.1
Wavelength (Å)	1.12	1.12
Space group	P2 ₁ 2 ₁ 2 ₁	P2 ₁ 22 ₁
Unit cell parameters (Å, °)	a = 43.8, b = 56.2, c = 66.4 Å α = β = γ = 90.0°	a = 47.3, b = 51.1, c = 72.9 Å α = β = γ = 90.0°
Resolution (Å)	50 - 1.55 (1.59-1.55) ^a	50 - 1.25 (1.28-1.25) ^a
Observations	172,015	334,858
Unique Reflections	24,187 (1,729) ^a	47,733 (2,877) ^a
Redundancy	7.1 (7.2) ^a	7.0 (5.7) ^a
Completeness (%)	98.7 (97.6) ^a	96.0 (79.6) ^a
<I/σ _I >	16.6 (1.0) ^a	18.2 (1.2) ^a
CC1/2	0.99 (0.42)	1.00 (0.65)
R _{sym} ^b	0.07 (1.98) ^a	0.05 (1.53) ^a
Refinement statistics		
Resolution (Å)	50 – 1.55	50 – 1.25
Reflections (work)	22,937	44,703
Reflections (test)	1,207	1,758
R _{cryst} (%) ^c / R _{free} (%) ^d	19.1 / 22.0	17.5 / 19.4
Average B (Å ²)	35.1	26.6
Wilson B (Å ²)	22.5	13.8
Protein atoms	3,667	3725
Waters	91	191
Other	0	0
RMSD from ideal geometry		
Bond length (Å)	0.005	0.005
Bond angles (°)	0.77	0.84
Ramachandran statistics (%) ^e		
Favored	99.1	99.5
Outliers	0.0	0.0
PDB Code	5CWB	5CWC

Supplementary Table 7. Data collection and refinement statistics for DHR7 and DHR8

Data collection	DHR7	DHR8
Beamline	ALS 8.3.1	ALS 8.3.1
Wavelength (Å)	1.12	1.12
Space group	C2	P2 ₁ 2 ₁ 2
Unit cell parameters (Å, °)	a = 84.7, b = 29.3, c = 63.2 Å α = γ = 90.0°, β = 123.5	a = 75.7, b = 76.4, c = 33.1 Å α = β = γ = 90.0°
Resolution (Å)	50 - 2.60 (2.67-2.60) ^a	50 - 1.80 (1.85-1.80) ^a
Observations	13,960	512,567
Unique Reflections	4,124 (322) ^a	72,203 (5,260) ^a
Redundancy	3.4 (3.4) ^a	7.1 (7.2) ^a
Completeness (%)	98.3 (97.9) ^a	99.9 (100.0) ^a
<I/σ _I >	4.9 (1.2) ^a	12.1 (1.1) ^a
CC1/2	0.99 (0.63)	1.00 (0.48)
R _{sym} ^b	0.15 (0.81) ^a	0.09 (1.85) ^a
Refinement statistics		
Resolution (Å)	50 – 2.60	50 – 1.80
Reflections (work)	3,916	68,562
Refections (test)	206	3,591
R _{cryst} (%) ^c / R _{free} (%) ^d	25.3 / 30.9	20.5 / 22.8
Average B (Å ²)	66.5	44.8
Wilson B (Å ²)	55.6	27.7
Protein atoms	2,521	11,104
Waters	0	466
Other	0	31
1GRMSD from ideal geometry		
Bond length (Å)	0.004	0.005
Bond angles (°)	0.58	0.73
Ramachandran statistics (%) ^e		
Favored	99.4	99.3
Outliers	0.0	0.0
PDB Code	5CWD	5CWF

Supplementary Table 8. Data collection and refinement statistics for DHR10 and DHR14

Data collection	DHR10	DHR14
Beamline	ALS 8.3.1	ALS 8.3.1
Wavelength (Å)	1.12	1.12
Space group	P2 ₁ 22 ₁	C222 ₁
Unit cell parameters (Å, °)	a = 47.6, b = 50.6, c = 72.4 Å α = β = γ = 90.0°	a = 41.0, b = 104.4, c = 90.2 Å α = β = γ = 90.0°
Resolution (Å)	50 - 1.20 (1.23-1.20) ^a	50 - 1.30 (1.33-1.30) ^a
Observations	313,175	964,864
Unique Reflections	54,627 (3,516) ^a	47,853 (3,495) ^a
Redundancy	5.6 (2.8) ^a	20.2 (17.0) ^a
Completeness (%)	98.5 (87.4) ^a	99.7 (99.8) ^a
<I/σ _I >	22.7 (1.9) ^a	19.2 (1.5) ^a
CC1/2	1.00 (0.81)	1.00 (0.85)
R _{sym} ^b	0.04 (0.53) ^a	0.08 (2.51) ^a
Refinement statistics		
Resolution (Å)	50 – 1.20	50 – 1.30
Reflections (work)	51,813	45,170
Refections (test)	2,730	2,377
R _{cryst} (%) ^c / R _{free} (%) ^d	14.1 / 16.9	16.2 / 18.6
Average B (Å ²)	21.3	30.1
Wilson B (Å ²)	12.0	15.2
Protein atoms	3,839	2,860
Waters	238	240
Other	10	20
1GRMSD from ideal geometry		
Bond length (Å)	0.007	0.018
Bond angles (°)	1.07	1.52
Ramachandran statistics (%) ^e		
Favored	99.6	100.0
Outliers	0.0	0.0
PDB Code	5CWG	5CWH

Supplementary Table 9. Data collection and refinement statistics for DHR18 and DHR49

Data collection	DHR18	DHR49
Beamline	ALS 8.3.1	ALS 8.3.1
Wavelength (Å)	1.12	1.12
Space group	P1	P1
Unit cell parameters (Å, °)	$a = 24.9, b = 43.7, c = 48.8$ $\alpha = 116.6, \beta = 90.7, \gamma = 102.2$	$a = 46.7, b = 61.7, c = 61.7$ $\alpha = 81.9, \beta = 67.8, \gamma = 67.9$
Resolution (Å)	50 – 1.75 (1.80-1.75) ^a	50 - 1.80 (1.85-1.80) ^a
Observations	36,231	103,465
Unique Reflections	17,261 (1,247) ^a	52,418 (3,824) ^a
Redundancy	2.1 (2.1) ^a	2.0 (2.0) ^a
Completeness (%)	94.8 (94.3) ^a	95.6 (94.5) ^a
$\langle I/\sigma_I \rangle$	9.8 (1.1) ^a	9.9 (1.4) ^a
CC1/2	1.00 (0.54)	1.00 (0.57)
$R_{\text{sym}}^{\text{b}}$	0.05 (0.85) ^a	0.04 (0.59) ^a
Refinement statistics		
Resolution (Å)	50 – 1.75	50 – 1.80
Reflections (work)	16,378	47,180
Refections (test)	866	1,864
$R_{\text{cryst}}(\%)^{\text{c}} / R_{\text{free}}(\%)^{\text{d}}$	18.5 / 22.2	24.0 / 28.8
Average B (Å ²)	41.2	49.4
Wilson B (Å ²)	23.9	28.9
Protein atoms	3,757	10,044
Waters	69	370
Other	5	0
1GRMSD from ideal geometry		
Bond length (Å)	0.010	0.005
Bond angles (°)	1.01	0.72
Ramachandran statistics (%) ^e		
Favored	99.2	98.7
Outliers	0.4	0.0
PDB Code	5CWI	5CWJ

Supplementary Table 10. Data collection and refinement statistics for DHR53 and DHR54

Data collection	DHR53	DHR54
Beamline	ALS 8.3.1	ALS 8.3.1
Wavelength (Å)	1.12	1.12
Space group	C2	P2 ₁ 2 ₁ 2 ₁
Unit cell parameters (Å, °)	a = 127.4, b = 24.4, c = 60.5 α = γ = 90.0, β = 116.9	a = 55.3, b = 69.1, c = 82.6 α = β = γ = 90.0
Resolution (Å)	50 – 1.90 (1.95-1.90) ^a	50 - 1.50 (1.54-1.50) ^a
Observations	46,741	360,741
Unique Reflections	13,327 (949) ^a	51,178 (3,720) ^a
Redundancy	3.5 (3.6) ^a	7.0 (7.1) ^a
Completeness (%)	98.3 (99.2) ^a	99.6 (99.9) ^a
<I/σ _I >	13.2 (1.0) ^a	15.1 (1.0) ^a
CC1/2	1.00 (0.74)	1.00 (0.49)
R _{sym} ^b	0.05 (1.67) ^a	0.06 (1.80) ^a
Refinement statistics		
Resolution (Å)	50 – 1.90	50 – 1.50
Reflections (work)	12,581	45,069
Refections (test)	663	1,817
R _{cryst} (%) ^c / R _{free} (%) ^d	24.5 / 28.5	18.3 / 22.2
Average B (Å ²)	83.9	43.1
Wilson B (Å ²)	38.0	23.2
Protein atoms	2,755	5,515
Waters	46	168
Other	0	0
1GRMSD from ideal geometry		
Bond length (Å)	0.006	0.011
Bond angles (°)	0.81	1.25
Ramachandran statistics (%) ^e		
Favored	97.7	98.9
Outliers	0.6	0.3
PDB Code	5CWK	5CWL

Supplementary Table 11. Data collection and refinement statistics for DHR64 and DHR71

Data collection	DHR64	DHR71
Beamline	ALS 8.3.1	ALS 8.3.1
Wavelength (Å)	1.12	1.12
Space group	P6 ₅	P2 ₁ 2 ₁ 2 ₁
Unit cell parameters (Å, °)	a = b = 98.3, c = 48.6 α = β = 90.0, γ = 120.0	a = 23.3, b = 82.9, c = 103.3 α = β = γ = 90.0
Resolution (Å)	50 - 2.90 (2.98-2.90) ^a	50 - 1.70 (1.74-1.70) ^a
Observations	67,571	158,634
Unique Reflections	6,086 (448) ^a	22,958 (1,666) ^a
Redundancy	11.1 (11.3) ^a	6.9 (7.1) ^a
Completeness (%)	100.0 (100.0) ^a	99.6 (100.0) ^a
<I/σ _I >	14.4 (1.3) ^a	15.0 (1.3) ^a
CC1/2	1.00 (0.44)	1.00 (0.64)
R _{sym} ^b	0.15 (2.14) ^a	0.08 (1.64) ^a
Refinement statistics		
Resolution (Å)	50 – 2.90	50 – 1.70
Reflections (work)	5,268	21,747
Refections (test)	277	1,144
R _{cryst} (%) ^c / R _{free} (%) ^d	21.4 / 26.5	20.6 / 25.0
Average B (Å ²)	104.9	36.1
Wilson B (Å ²)	72.9	24.4
Protein atoms	3,656	3,257
Waters	29	95
Other	0	0
1GRMSD from ideal geometry		
Bond length (Å)	0.019	0.010
Bond angles (°)	1.70	1.14
Ramachandran statistics (%) ^e		
Favored	93.3	99.0
Outliers	0.9	0.0
PDB Code	5CWM	5CWN

Supplementary Table 12. Data collection and refinement statistics for DHR76 and DHR79

Data collection	DHR76	DHR79
Beamline	ALS 8.3.1	ALS 8.3.1
Wavelength (Å)	1.12	1.12
Space group	P2 ₁ 2 ₁ 2 ₁	P4 ₃ 2 ₁ 2
Unit cell parameters (Å, °)	a = 62.5, b = 66.3, c = 102.2 α = β = γ = 90.0°	a = b = 94.3, c = 77.6 Å α = β = γ = 90.0°
Resolution (Å)	50 – 3.55 (3.44-3.55) ^a	50 - 1.90 (1.95-1.90) ^a
Observations	45,810	398,338
Unique Reflections	6,514 (470) ^a	28,192 (2,042) ^a
Redundancy	7.0 (7.2) ^a	14.1 (14.4) ^a
Completeness (%)	99.8 (100.0) ^a	100.0 (100.0) ^a
<I/σ _I >	10.9 (1.2) ^a	17.5 (1.1) ^a
CC1/2	1.00 (0.33)	1.00 (0.48)
R _{sym} ^b	0.18 (1.91) ^a	0.11 (3.61) ^a
Refinement statistics		
Resolution (Å)	50 – 3.55	50 – 1.90
Reflections (work)	6,162	26,730
Refections (test)	322	1,406
R _{cryst} (%) ^c / R _{free} (%) ^d	37.2 / 40.8	19.3 / 23.8
Average B (Å ²)	242.9	58.6
Wilson B (Å ²)	111.5	34.2
Protein atoms	6,992	3,591
Waters	0	63
Other	0	0
1GRMSD from ideal geometry		
Bond length (Å)	0.004	0.013
Bond angles (°)	0.83	1.29
Ramachandran statistics (%) ^e		
Favored	98.6	98.3
Outliers	0.0	0.4
PDB Code	5CWO	5CWP

Supplementary Table 13. Data collection and refinement statistics for DHR81

Data collection	DHR81
Beamline	ALS 8.2.1
Wavelength (Å)	1.00
Space group	I23
Unit cell parameters (Å, °)	$a = b = c = 120.0$ $\alpha = \beta = \gamma = 90.0$
Resolution (Å)	50 - 2.05 (2.10-2.05) ^a
Observations	397,887
Unique Reflections	18,193 (1,317) ^a
Redundancy	21.9 (22.4) ^a
Completeness (%)	100.0 (100.0) ^a
$\langle I/\sigma_I \rangle$	16.8 (1.1) ^a
CC1/2	1.00 (0.32)
$R_{\text{sym}}^{\text{b}}$	0.31 (3.53) ^a
Refinement statistics	
Resolution (Å)	50 – 2.05
Reflections (work)	16,326
Reflections (test)	860
$R_{\text{cryst}}(\%)^{\text{c}} / R_{\text{free}}(\%)^{\text{d}}$	21.7 / 26.6
Average B (Å ²)	49.3
Wilson B (Å ²)	31.4
Protein atoms	3,654
Waters	51
Other	14
1GRMSD from ideal geometry	
Bond length (Å)	0.007
Bond angles (°)	0.85
Ramachandran statistics (%) ^e	
Favored	98.3
Outliers	0.4
PDB Code	5CWQ

Supplementary Table 14. Data collection and phasing statistics for DHR5 (S-SAD)

Data collection

Beamline	ALS 8.3.1
Space group	P2 ₁ 22 ₁
Unit cell parameters (Å, °)	a = 47.3, b = 50.7, c = 73.0 α = β = γ = 90.0
Wavelength (Å)	1.714
Resolution (Å)	50 – 1.85 (1.90-1.85) ^a
Observations	556,066
Unique Reflections	28,741 (1,966) ^a
Redundancy	19.3 (7.2) ^a
Completeness (%)	99.2 (91.0) ^a
<I/σ _I >	31.9 (3.2) ^a
CC1/2	1.00 (0.92)
R _{sym} ^b	0.08 (0.50) ^a
FOM	0.58
PDB Code	5CWC

Supplementary Tables 6-14 | Crystallography data collection

^a Numbers in parentheses refer to the highest resolution shell.

^b $R_{\text{sym}} = \sum_{hkl} \sum_i |I_{hkl,i} - \langle I_{hkl} \rangle| / \sum_{hkl} \sum_i I_{hkl,I}$, where $I_{hkl,i}$ is the scaled intensity of the i^{th} measurement of reflection h, k, l, $\langle I_{hkl} \rangle$ is the average intensity for that reflection, and n is the redundancy ²⁴.

^c $R_{\text{cryst}} = \sum_{hkl} |F_o - F_c| / \sum_{hkl} |F_o| \times 100$

^d R_{free} was calculated as for R_{cryst} , but on a test set comprising 5% of the data excluded from refinement.

^e Calculated using Molprobity²⁵.

Supplementary Discussion 5 | Small Angle X-ray Scattering (SAXS) analysis

Guinier and P(r) analysis were done using ATSAS²⁶. The Porod exponent was determined from a linear regression analysis (I vs q) of the top of the first peak in the Porod-Debye plot ($q^4 \cdot I(q)$ vs q^4) of the scattering data, implemented in SCATTER, available at beamline 12.3.1^{27,28}. The molecular mass in solution was calculated using SCATTER²⁹. The results are summarized in Supplementary Table 15.

25% of the designs had molecular weights in solution that were significantly greater than the predicted molecular weight (1.2-4 fold), suggesting that these designs formed multimeric assemblies or a small portion of aggregates²⁹. All 55 designs had Porod exponents (P_E) greater than 2.9, indicating significant levels of folded protein; 67% of the designs had a P_E of 3.4-4, indicating a well-folded core²⁸. Of the 15 proteins that crystallized, the majority (66%) had P_E of 3.9-4, consistent with more well-packed proteins being easier to crystallize.

Radius of gyration (R_g) and maximum of distance distribution (d_{max}) were calculated from real space distance distribution $P(r)$. Distance distribution comparisons between experimental data and models are included in supporting_experimental_data.pdf. Among the models confirmed by crystallography, DHR 49 and 76 formed dimers in solution. The experimental data were fit using models based on the dimer configuration observed in the crystal structure. DHR 5 tendency to aggregation (see SEC in supporting_experimental_data.pdf) affected the SAXS profile resulting in a high Molecular weight and V_r above our acceptance threshold.

If molecular mass and R_g of models were within a 25% error from experimental data and V_r was below 2.5, the models were considered able to recapture the SAXS data. D_{max} errors are generally within 25%; higher values are highlighted in Supplementary Table 15.

43 designs satisfied our requirements: DHR 1 2 3 4 7 8 9 10 14 15 18 20 21 23 24 26 27 31 32 36 39 46 47 49 52 53 54 55 57 58 59 62 64 68 70 71 72 76 77 78 79 80 81 82. The corresponding models, including also DHR5, are shown in Extended Data Fig. 9.

	EKSTDEEEIRELLQRAEERIREAQERCREGDGWLLEHHHHHH
DHR27	MTRQEOLDEVLEBTORLAEARKLMTDEEEAKKIQEEAERAKEMLRRAVEKVTNDNEVIEKLLEVVKETIRLAEEAMKKMTDEEEAAKI AKEALEAIKMLARAVEEVTDNEVIEKLLEVVKETIRLAEEAMKKMTDEEEAAKI AKEALEAIKMLARAVEEVTDKERIEQLLREVKEEI RRAEEESRKETDDEEAAGRAREALRRIRERAREVEEDKSGWLEHHHHHH
DHR28	MDEEVQRIEEVRRAEVRESLERNDSEEAEELAREALERVAEEVKESIKERPDRLIAIEAIRALVRLAIEIVRLALEQNDSELAREV AEEALRAVAEVVKEAIRQGRDRDLIAIEAIRALVRLAIEIVRLALEQNDSELAREVAAEALRAVAEVVKEAIRQGRDRELAKAIRALR LAEEIRRLAEEQNDDELAREVEELAREAIEEVRKELERQRPGRGWLEHHHHHH
DHR29	MSEVEESAQEVEKRAQEVRREEAERRGTSQEVLDEIKRVVDEARQLAQAKESDDSEVAESALQVVREALKVVLALARGTSEEVLK RVVSEAIKLAIAKSSDSEVAESALQVVREALKVVLALARGTSEEVLK LEQLERGTSEELRESREVSENIRKALEEIKSPDGWLEHHHHHH
DHR30	MSTVKELLDRARELRLAERASFQGSDEEEAKRLLEDLEQLVQEIRRELEETGTSSEVIRLIAKAIMLMAELALRAAEQGSDAAEAMK LLKDLLRLVLEILRELRETGTDSEVIRLIAKAIMLMAELALRAAEQGSDAAEAMKLLKDLLRLVLEILRELRETGTDKEIRKVAEEIM RRAKTALDEARQGSDAAEAMKRLKEQRLRILERLREEREKGTDGWLEHHHHHH
DHR31	MDSYTERARKAVKRYVKEEGGSFFEEAEREAEVKVREEIRKKASDSYLIQAAAAAVVAYVIEEGGSPEEAVKIAEEVRRRIKEADDSYLIQ AAAAVVAYVIEEGGSPEEAVKIAEEVRRRIKEADDRELIRRAEAEVIERGGSPEEAVKEAEKEVKKQKEESDGWLEHHHHHH
DHR32	MSIQEKAKQSVKRVEEGGSSEEAEERAKERAKEEVLKKEADDSTLVRAAAAVVLYVLEKGGSTEEAVQAREVIERLKEASDSTLVRA AAAUVLYVLEKGGSTEEAVQAREVIERLKEASDEELIREAAKEVKKLKEEGGSVVEAVERARERERIEELQKRSDDGWLEHHHHHH
DHR33	MSETTEVKKLVEEKVKKEGGSPPEEAKETAKEVTEELKEESQDSTLLKVALVASAVLKKEGGSPPEEAAETAKEVVKELRKSASDSTLLKV AALVASAVLKKEGGSPPEEAAETAKEVVKELRKSASDEELLKEAARQAEESLRQGKSPPEEAAEAKKEVKKLKEKSQDGWLEHHHHHH
DHR34	MSETTEVKKLCEEKVKKEGGSPPEEAKETAKEVTEELKEESQDSTLLKVALCASAVLKKEGGSCEEAAETAKEVVKELRKSASDSTLLKV AALCASAVLKKEGGSPPEEAAETAKEVVKELRKSASDEELLKEAARQAEESLRQGKSCFEEAAEAKKEVKKLKEKSQDGWLEHHHHHH
DHR35	MSEDEVAQQSRYAKEQGGDPEKSREEAEKALEEVVKQATSSEALQVALEAARYASEEGEDPAEALKEAARALEEVRRSATSSSEALQV ALEAARYASEEGEDPAEALKEAARALEEVRSATSEEDILKEALDRAREASERGQNPNAELKEAAEELKKKEKSDGWLEHHHHHH
DHR36	MSDLEKAKRKFVKEEKKGRNPEEAKKLKLLKKSAGSSDLTALAKFVLEEVKGRNPKRVAEEAIKQAKEDRKRNSNGWLEHHHHHH LAKFVLEEVKGRNPEEAVKEAIKLAELKLRKSAGSSEQELEKLATKVLLEEVKKGRNPKRVAEEAIKQAKEDRKRNSNGWLEHHHHHH
DHR37	MSSTERAAQSVKYLQQQGKDQAKKAQEVKENIEKEANSSSVIRAAAAVVFYLLEQGYDPDQALKKAQEVARNIENEANSSSVIRA AAAUVFYLLEQGYDPDQALKKAQEVARNIENEANSDDVIKEAKVVKRLEEGQDPDKALEEARKRAQKTEKTTSGWLEHHHHHH
DHR38	MSSTERAAQSCKKYLQQQGKDQAKKAQEVKENIEKEANSSSVIRAAAACVFYLLEQGYDCDQALKKAQEVARNIENEANSSSVIRA AAAUVFYLLEQGYDCDQALKKAQEVARNIENEANSDDVIKEAKVVKRLEEGQDCDKALEEARKRAQKTEKTTSGWLEHHHHHH
DHR39	MSDLQEVADRIEQLKREGRSPPEEARKEARLIEEIKQSAGGDSELIEAVRIVKELEEQGRSPSEAACEVELIERIRRAAGGDSELI EVAVRIVKELEEQGRSPSEAACEVELIERIRRAAGGDSDRIKKAVELVRELEERGRSPSEAARRAVEEIQRSVEEDGGNGWLEHHHHH H
DHR40	MSESDEVAKRISKEAKKEGRSSEEVVELVERFREAIEKLKEQGDSEAIRVAVEIAADEALREGLSPEEVVELVERFVQAIQKLQENGSE AIRVAVEIAADEALREGLSPEEVVELVERFVQAIQKLQENGEEDEIQKAVETAQEQLEERGRSPKEVETVETVEEQVKEVEEKQKGEGWLEH HHHHH
DHR41	MSDIKEKAKRIADRAIDVVRKAAEKEGGSPPEKIREALQQAKRCAEKLIRLVKEAQESNSSDVREAARVVALEAVRVVVRAAEKGGSPPEV VEAVCRAVRCAEKLIRLVKRAEESNSSDVREAARVVALEAVRVVVRAAEKGGSPPEEVVEAVCRAVRCAEKLIRLVKRAEESNSENVRES ARRAELKVLKTVQQAEEEGKSPPEEVVEQCRSVRKAEEQI RETQERERSTSGWLEHHHHHH
DHR42	MSDAEVKKQAEFIANRAYTAQKQGESDSRAKKAELKVRKAAEKLARIERAQKEGDSDALEVARQALEIARRAFETAKKQGHSATEA AKAFVDVVEAAISLAELIISAKRQGDSDALEVARQALEIARRAFETAKKQGHSATEAAKAFVDVVEAAISLAELIISAKRQGDQKALEI ARKALQAKAKENFEEAQKRGESATQAQAKRFVDTVEKEIKKAQEQIKRERKGDWLEHHHHHH
DHR43	MSKEELIEKARRVAKEAIEEAKRQGDPSEAKKAAEKLIAKVEEAVKEAKRLKEEGNSELAELISEAIQVAVEAVEAVRQGKDPFKA AEAAEELIRAVVVAEAVKEAERLKREGNSELAELISEAIQVAVEAVEAVRQGKDPFKAEEAAEELIRAVVVAEAVKEAERLKREGNSELA AKINDTIREAVREVQQAEEGKDPFEEAREAAEAKIRESERVREEEEKKRNGWLEHHHHHH
DHR44	MSNEQEKKDLKKAEEAKSPDPFEEAREAAEAKIRESERVREEEEKKRNGWLEHHHHHH PDPELIRLAIEAAERSGSEKAEIILRAEEAQSPPDPELQKLAKERLGGWLEHHHHHH
DHR45	MSSEEELEKDAREASESGADPEWLREIVDLARESGDSEVIELAKRALEAAKSGADPEWLLRIVRQAEESGSSEVIELAKRALEAKSG ADPEWLLRIVRQAEESGSSEEVIELAKRALEAAKKGKDPKELLEEVKRKEESGGWLEHHHHHH
DHR46	MSTKEEKERIERIEKEVRSPDPENIREAVRKAEEELLRENTPSTEAEELLRAIAEAVRAPDPEAIREAVRAAEELLRENTPSTEAEELLRR AIEAAVRAPDPEAIREAVRAAEELLRENTPSEEAKLLRRAIESAKAPDPEAQREAKRAEEELRKEDPGWLEHHHHHH
DHR47	MSTKEEKERIERIEKEVRSPDCENIREAVRKAEEELLRENTPSTEAEELLRAIAEAVRCPDCEAIREAVRAAEELLRENTPSTEAEELLRR AIEAAVRCPDCEAIREAVRAAEELLRENTPSEEAKLLRRAIESAKCKPDPPEAQREAKRAEEELRKEDPGWLEHHHHHH
DHR48	MNSREEEEAKRIVKEAKKSGFDPEEVKEKALREVIRVAEEGTGSEALKEALKIVEEAKSGYDPAEVAKALAEVIRVAEEGTGSEALKEA LKIVEEAAKSGYDPAEVAKALAEVIRVAEEGTGSEALKEALKIVEEAKSGYDPAEVAKALAEVIRVAEEGTGSEALKEA
DHR49	MDSEEQERIRRILKEARKSGTEESLRQAIEDVAQLAKKSQDSEVLEEAIRVILRIAKESGSEEALRQAIRVAETAKEAQDSEVLEEA IRVILRIAKESGSEEALRQAIRVAETAKEAQDPRVLEEAIRVIRQIAEESGSEEARRQAEREEEIRRAQGWLEHHHHHH
DHR50	MDPEVVRREVERATEEYRKNPGSDEAREOLKEAVERAEEAARSPDPEAVQAVAVEAATQIYENTPGSEEAKKALEIAVRAAENAARLPDP EAVQAVAVEAATQIYENTPGSEEAKKALEIAVRAAENAARLPDPPEAVRVAEEAADQIRKNTPGSELAKRADEIKKRARELLERLPGWLEH HHHHH
DHR51	MQSEDRKEKIRELERKARENTGSDEARQAVKEIARIAKEALEEGNADTAKEAIQRLEDLARDYSGSDVASLAVKAIKIAETALRNGYA DTAKEAIQRLEDLARDYSGSDVASLAVKAIKIAETALRNGYKETAEEAIKRLRELAEDYKGSEVAKLAAEAIERIEKVSRSERGGWLEH HHHHH
DHR52	MQCDRKEKIRELERKARENTGSDEARQAVKEIARIAKEALEEGCCDTAKEAIQRLEDLARDYSGSDVASLAVKAIKIAETALRNGCC DTAKEAIQRLEDLARDYSGSDVASLAVKAIKIAETALRNGCKETAEEAIKRLRELAEDYKGSEVAKLAAEAIERIEKVSRSERGGWLEH HHHHH
DHR53	MSNDEKEKLKELLKRAEELAKSPDPDELIKEAVRLAEEVVRERPGSNLAKKALEIIILRAAEELAKLPDPEALKEAVKAAEVVREOPGSN LAKKALEIIILRAAEELAKLPDPEALKEAVKAAEVVREOPGSNLLAKKALEIIILRAAEELKKSPDPEAQKBAKKAQKVREERPGGWLEH HHHHH
DHR54	MTTEDERRELEKVARAKIAEAREGNTDEVREQLQRALEIARESGTTEAVKLALEVVARVATEAARRGNTDAVREALEVALEIARESGTT EAVKLALEVVARVATEAARRGNTDAVREALEVALEIARESGTEEAVRLALEVVKRVSDEAKKQGNEDAVERKKIEEESGGWLEH HHHHH

	KIVKAIQEAVESLREAAEESGDPEKREKARERVREAVERAEEVQRDPGWLEHHHHHH
DHR80	MNSEELEREEAERRLOEARKRSEEARERGDLKELAEALIEEARAVOELARVASERGNSEEAERASEKAQRVLEEARKVSEEAREQGD DEVLALALIAIALAVLALAEVASSRGNSEEAERASEKAQRVLEEARVKSEEAREQGDDEVLALALIAIALAVLALAEVASSRGNKEEAE RAYEDARRVEEEARKVKESAEEQGDSEVKRLAEEAEQLAREARRHVQETRGGWLEHHHHHH
DHR81	MNSEELEREEAERRLOEARKRSEEARERGDLKELAEALIEEARAVOELARVACERGNSEEAERASEKAQRVLEEARVKSEEAREQGD DEVLALALIAIALAVLALAEVACCRGNSEEAERASEKAQRVLEEARVKSEEAREQGDDEVLALALIAIALAVLALAEVACCRGNKEEAE RAYEDARRVEEEARKVKESAEEQGDSEVKRLAEEAEQLAREARRHVQECRGGWLEHHHHHH
DHR82	MNDEEVQEAVERAELREEAEELIKKARKTGDPELLRKALEALEEAVRAVEEAIKRNPNDNEAVETAVRLARELKVAEELQERAKKTG DPELLKLALRALEVAVRAVELAIKSNPDNDECVTAVRLARELKVAEELQERAKKTGDPELLKLALRALEVAVRAVELAIKSNPDNEE AVETAKRLAELRKVAELLEERAKETGDPELQELAKRAKEVADRALAKKSNPNGWLEHHHHHH
DHR83	MNDEEVQEACERAELREEAEELIKKARKTGDPELLRKALEALEEAVRAVEEAIKRNPNDDECVTACRLARELKVAEELQERAKKTG DPELLKLALRALEVAVRAVELAIKSNPDNECVETACRLARELKVAEELQERAKKTGDPELLKLALRALEVAVRAVELAIKSNPDNEE CVETAKRLAELRKVAELLEERAKETGDPELQELAKRAKEVADRALAKKSNPNGWLEHHHHHH

Supplementary References

1. Leaver-Fay, A. *et al.* ROSETTA3: an object-oriented software suite for the simulation and design of macromolecules. *Methods Enzymol.* **487**, 545–574 (2011).
2. Rohl, C. A., Strauss, C. E. M., Misura, K. M. S. & Baker, D. in *Methods in Enzymology* (ed. Johnson, L. B. and M. L.) **383**, 66–93 (Academic Press, 2004).
3. Wintjens, R. T., Rooman, M. J. & Wodak, S. J. Automatic Classification and Analysis of $\alpha\alpha$ -Turn Motifs in Proteins. *J. Mol. Biol.* **255**, 235–253 (1996).
4. Koga, N. *et al.* Principles for designing ideal protein structures. *Nature* **491**, 222–227 (2012).
5. Kuhlman, B. & Baker, D. Native protein sequences are close to optimal for their structures. *Proc. Natl. Acad. Sci.* **97**, 10383–10388 (2000).
6. Sheffler, W. & Baker, D. RosettaHoles: Rapid assessment of protein core packing for structure prediction, refinement, design, and validation. *Protein Sci.* **18**, 229–239 (2009).
7. Parmeggiani, F. *et al.* A General Computational Approach for Repeat Protein Design. *J. Mol. Biol.* **427**, 563–575 (2015).
8. Kuhlman, B. *et al.* Design of a Novel Globular Protein Fold with Atomic-Level Accuracy. *Science* **302**, 1364–1368 (2003).
9. Jones, D. T. Protein secondary structure prediction based on position-specific scoring matrices. *J. Mol. Biol.* **292**, 195–202 (1999).
10. Bradley, P., Misura, K. M. S. & Baker, D. Toward High-Resolution *de Novo* Structure Prediction for Small Proteins. *Science* **309**, 1868–1871 (2005).
11. Das, R. *et al.* Structure prediction for CASP7 targets using extensive all-atom refinement with Rosetta@home. *Proteins Struct. Funct. Bioinforma.* **69**, 118–128 (2007).
12. Huang, P.-S. *et al.* RosettaRemodel: A Generalized Framework for Flexible Backbone Protein Design. *PLoS ONE* **6**, e24109 (2011).

13. Kajava, A. V. Tandem repeats in proteins: From sequence to structure. *J. Struct. Biol.* **179**, 279–288 (2012).
14. Di Domenico, T. *et al.* RepeatsDB: a database of tandem repeat protein structures. *Nucleic Acids Res.* **42**, D352–D357 (2014).
15. Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics* **22**, 1658–1659 (2006).
16. Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-generation sequencing data. *Bioinformatics* **28**, 3150–3152 (2012).
17. Zhang, Y. & Skolnick, J. TM-align: a protein structure alignment algorithm based on the TM-score. *Nucleic Acids Res.* **33**, 2302–2309 (2005).
18. Walsh, I. *et al.* RAPHAEL: recognition, periodicity and insertion assignment of solenoid protein structures. *Bioinformatics* **28**, 3257–3264 (2012).
19. Punta, M. *et al.* The Pfam protein families database. *Nucleic Acids Res.* **40**, D290–D301 (2012).
20. Kabsch, W. & Sander, C. Dictionary of protein secondary structure: Pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers* **22**, 2577–2637 (1983).
21. Kobe, B. & Kajava, A. V. When protein folding is simplified to protein coiling: the continuum of solenoid protein structures. *Trends Biochem. Sci.* **25**, 509–515 (2000).
22. Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.* **25**, 3389–3402 (1997).
23. Remmert, M., Biegert, A., Hauser, A. & Söding, J. HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment. *Nat. Methods* **9**, 173–175 (2012).

24. Weiss, M. S. & Hilgenfeld, R. On the use of the merging R factor as a quality indicator for X-ray data. *J. Appl. Crystallogr.* **30**, 203–205 (1997).
25. Chen, V. B. *et al.* MolProbity: all-atom structure validation for macromolecular crystallography. *Acta Crystallogr. Sect. D* **66**, 12–21 (2010).
26. Petoukhov, M. V. *et al.* New developments in the ATSAS program package for small-angle scattering data analysis. *J. Appl. Crystallogr.* **45**, 342–350 (2012).
27. Classen, S. *et al.* Implementation and performance of SIBYLS: a dual endstation small-angle X-ray scattering and macromolecular crystallography beamline at the Advanced Light Source. *J. Appl. Crystallogr.* **46**, 1–13 (2013).
28. Rambo, R. P. & Tainer, J. A. Characterizing flexible and intrinsically unstructured biological macromolecules by SAS using the Porod-Debye law. *Biopolymers* **95**, 559–571 (2011).
29. Rambo, R. P. & Tainer, J. A. Accurate assessment of mass, models and resolution by small-angle scattering. *Nature* **496**, 477–481 (2013).